

Sequence Diversity within STR Loci: Applications to Human Identification

Peter M. Vallone, Ph.D.

Leader, Applied Genetics Group

10th SFAF Meeting

Santa Fe, New Mexico

May 29, 2015



**National Institute of
Standards and Technology**

U.S. Department of Commerce

Disclaimer

This presentation will mention commercial products, but we are in no way attempting to endorse any specific products.

NIST Disclaimer: Certain commercial equipment, instruments and materials are identified in order to specify experimental procedures as completely as possible. In no case does such identification imply a recommendation or it imply that any of the materials, instruments or equipment identified are necessarily the best available for the purpose.

Information presented does not necessarily represent the official position of the National Institute of Standards and Technology or the U.S. Department of Justice.

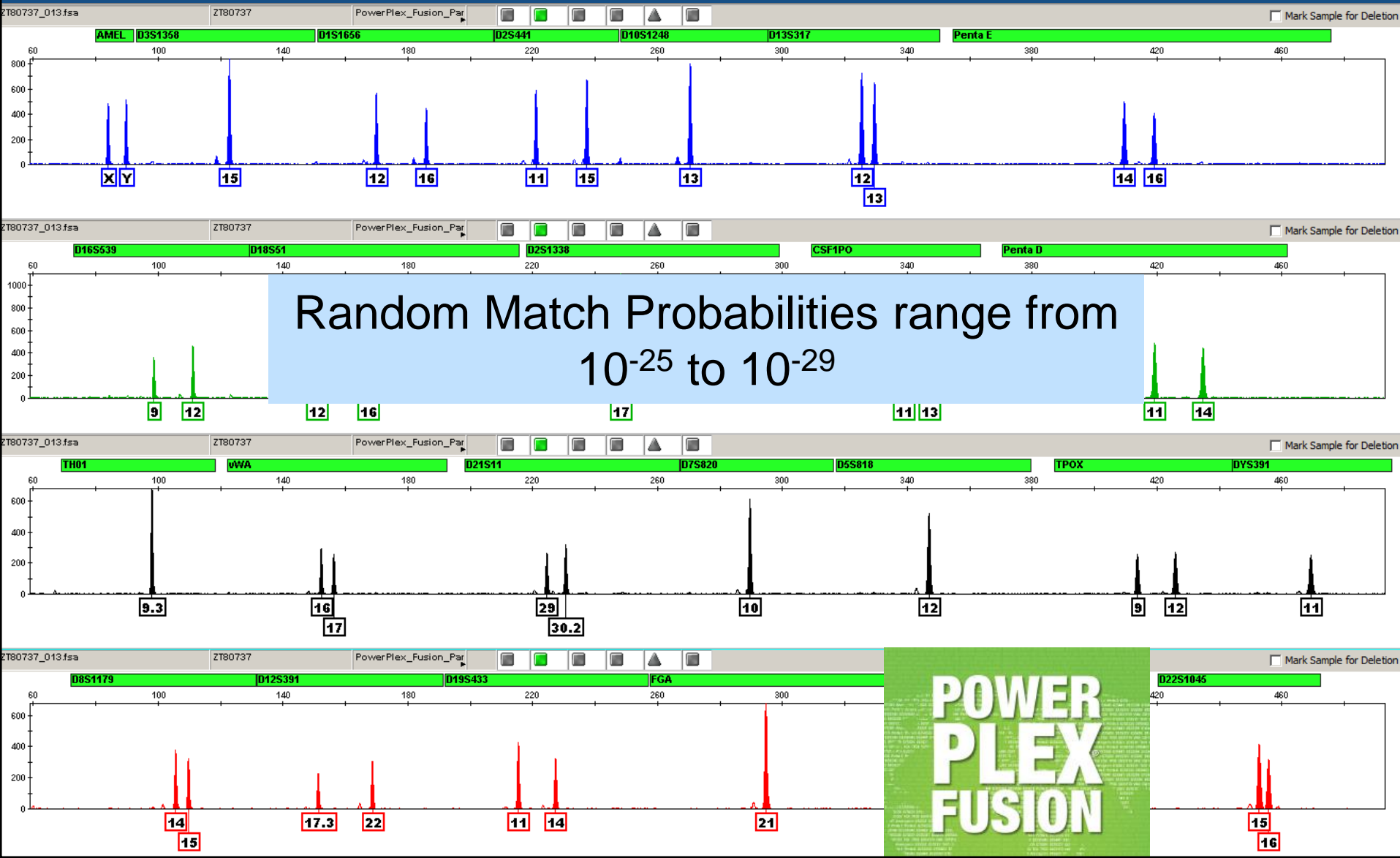
Our group receives or has received funding from the FBI Laboratory and the National Institute of Justice.

Starting with conclusions...

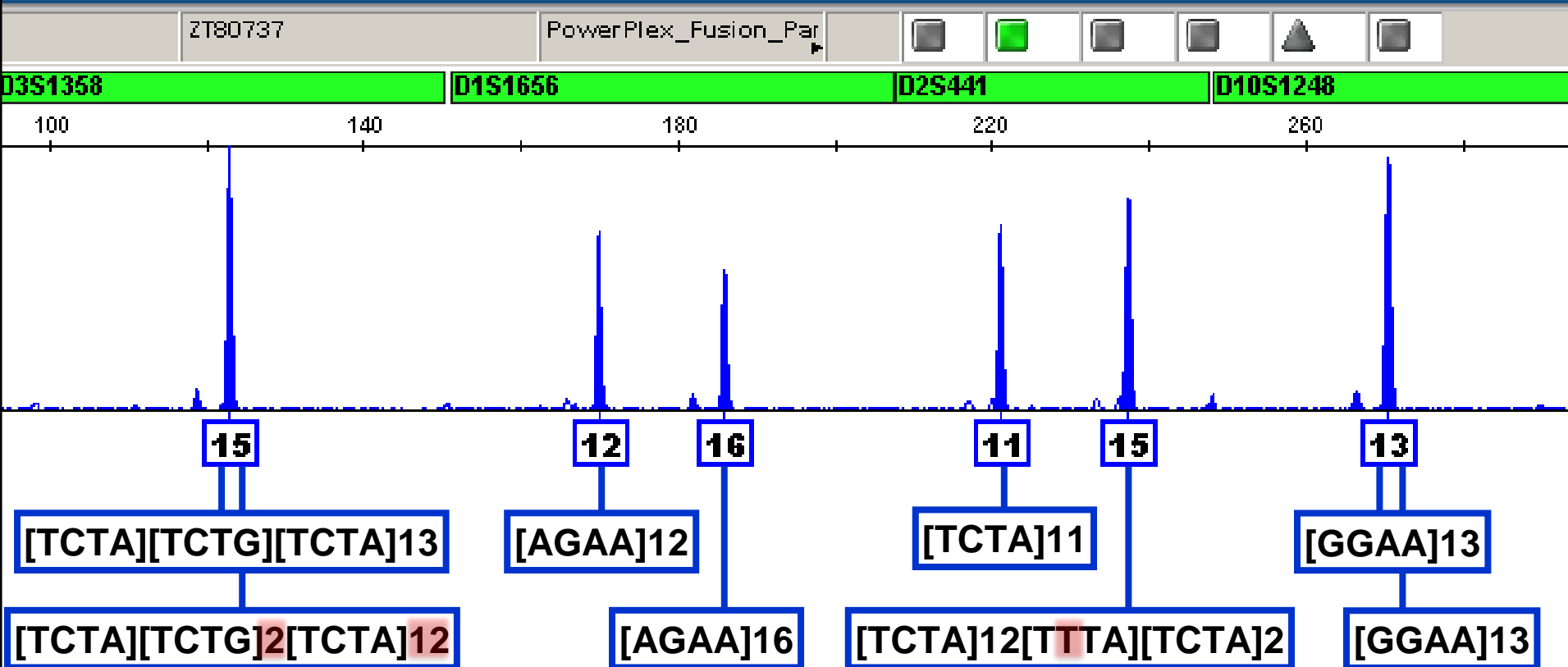
- 24 loci were sequenced for a set of 183 samples (22 autosomal STRs, 1 Y STR, and Amelogenin)
- Length and sequence variation within the STR region and surrounding flanking regions were detected by NGS and compared to length-based genotypes (CE)
- Loci with compound and complex repeat motifs contained the majority of 'additional' information by NGS
- Isoalleles and flanking region polymorphisms detected by NGS *may* assist in resolving mixtures...

Forensic STR Sequence Diversity

Length-based genotyping by CE methods



Forensic STR Sequence Diversity



Sequence-Based Heterozygote: A locus that appears homozygous in length-based measurements (such as CE), but is heterozygous by sequence

NGS has potential for finer resolution of STR amplicons not detectable by CE-length based methods

- Additional STR alleles
- Flanking region SNPs and InDels

Applications to DNA mixtures

- Resolve isoalleles (identical by length alleles)
- Resolve minor contributor peaks from stutter



*General gains of
larger multiplexing and SNP detection*

Forensic STR Sequence Diversity Methods

NIST Population Samples

- N=183
 - Caucasian (70)
 - Hispanic (45)
 - African American (68)

Amplification & Library Prep

- 1 ng input DNA
- PowerSeq Auto System (Promega)
- Illumina TruSeq HS PCR-Free

Sequencing

- MiSeq



Bioinformatics

- Sequence
- CE concordance



Pop Gen

- P_1 / HET
- Repeat Region
- Flank Analysis
- Stutter

Recognition Site-Based Informatics for STRs



Computer Returns:

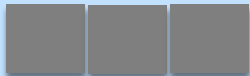
The length between the recognition sequences = 36 bp

Reference table in software returns a “9” allele

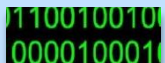
The sequence between the recognition sites



PCR Primers



STR Repeat Region



Recognition site (~10 nt)

¹ <http://battelleexactid.org/>

² STRait Razor: a length-based forensic STR allele-calling tool for use with second generation sequencing data. Warshauer et al., Forensic Sci Int Genet. 2013 (7):409-17

³ STRait Razor v2.0: the improved STR Allele Identification Tool--Razor. Warshauer et al., Forensic Sci Int Genet. 2015 (14):182-6

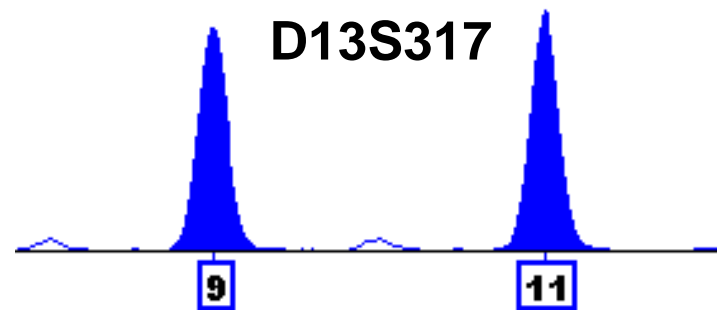
Forensic STR Sequence Diversity

CE (length-based genotype) concordance check results:

24 loci x 183 samples = 4392 loci evaluated

> 99% concordance with CE data

Forensic STR Sequence Diversity



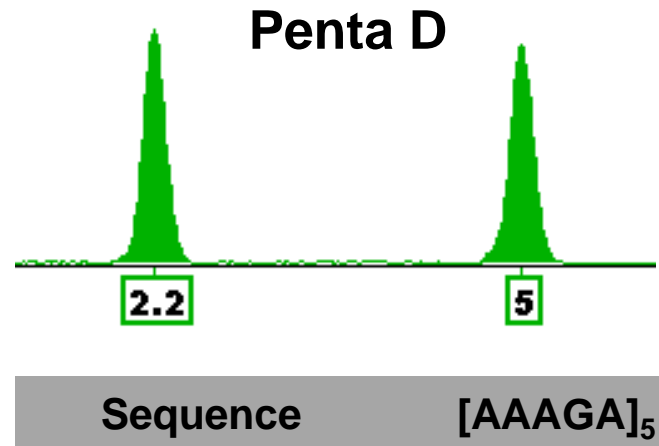
Sequence [TATC]₁₀, [TATC]₁₁

← Repeat Region NGS Recognition Region 4 bp Deletion CE Primer Binding Site→

TATC TATC TATC AATCAATCATCTATCTATCTTTCTGTC-----TTTTTGGGCTGCCTATGGCTCAA
TATC TATC TATC AATCAATCATCTATCTATCTTTCTGTCGTCTTTTTGGGCTGCCTATGGCTCAA

Flanking region InDel: Bioinformatic pipelines may reduce the region used for genotyping, resulting in deletions not being “counted” as they would via CE

Forensic STR Sequence Diversity



13 bp deletion in
Recognition Region

←CE Primer Binding Site

TAGGTTACAGAGCAAGACACCATCTCAAG-----AAAGA AAAGA AAAGA AAAGA AAA

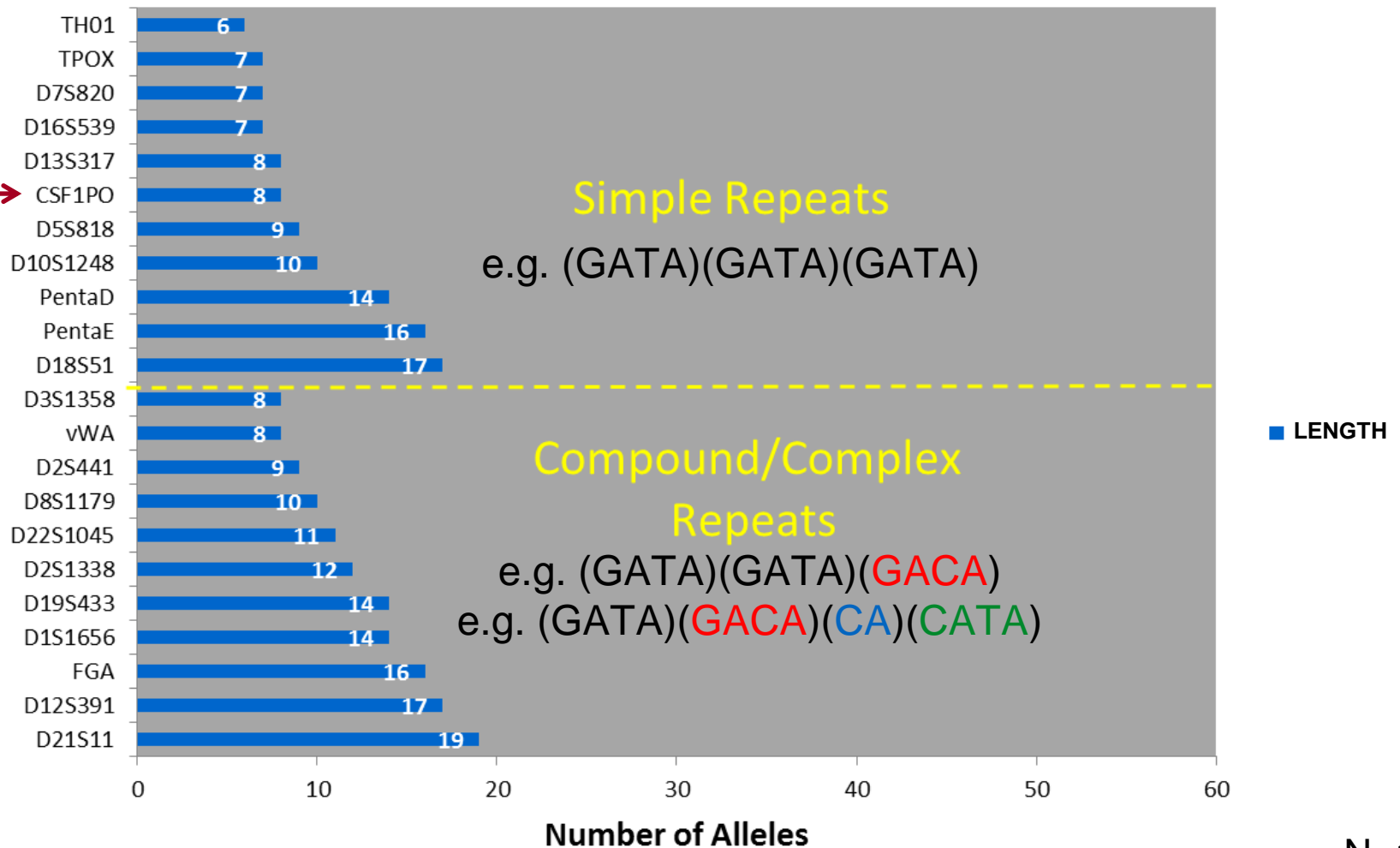
TAGGTGACAGAGCAAGACACCATCTCAAGAAAGAAAAAAAAAGAAAGA AAAGA AAAGA AAAGA AAA

Repeat Region →

Bioinformatic Null Allele: A true allele that is present within the raw sequence data but is not detected by the bioinformatic pipeline

Forensic STR Sequence Diversity

Alleles Obtained by Length



Forensic STR Sequence Diversity

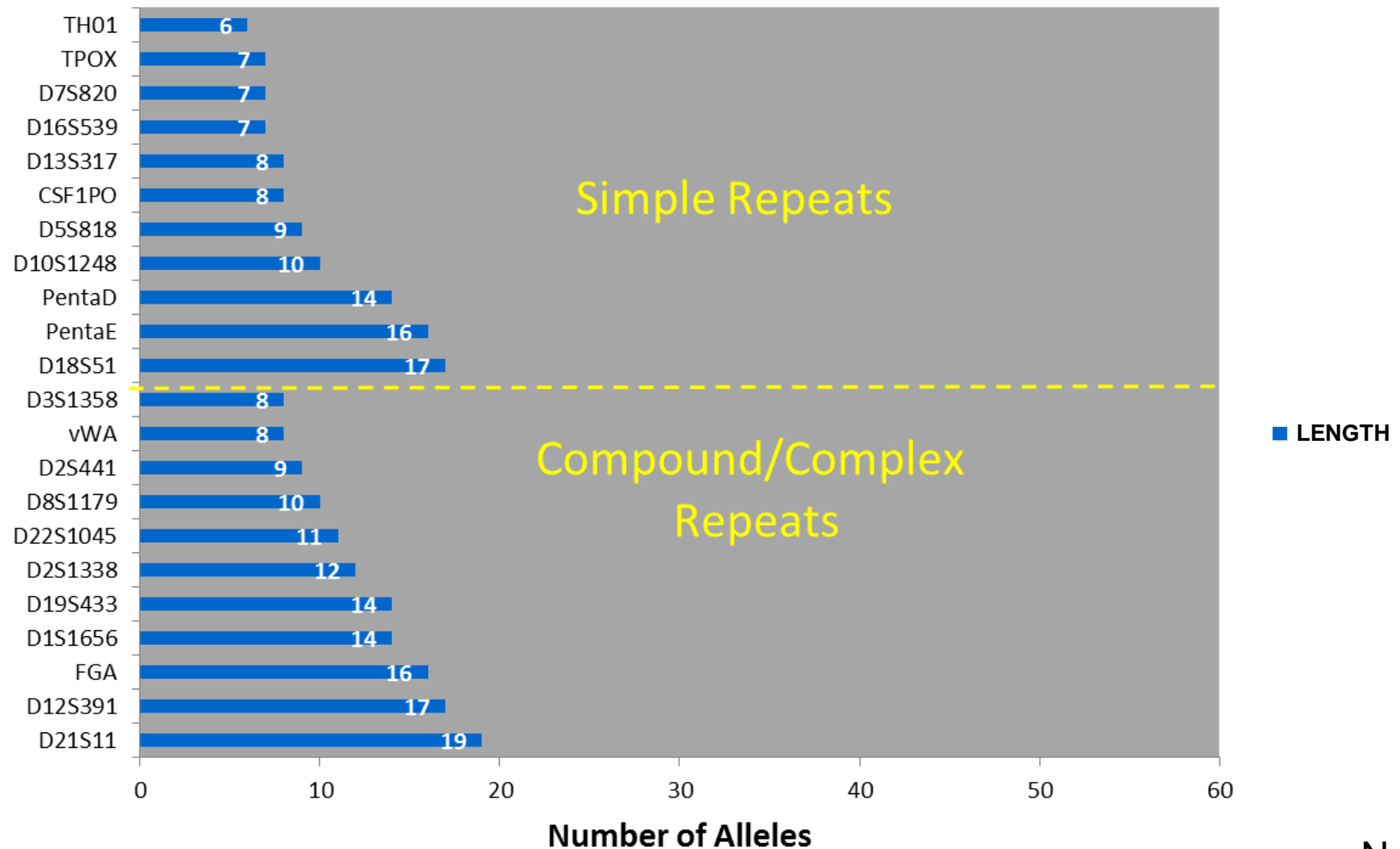
Additional Alleles by Sequence

CSF1PO														
7	[AGAT]7	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT						
8	[AGAT]8	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT					
9	[AGAT]9	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT				
10	[AGAT]10	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT			
10	[AGGT][AGAT]9	AG	G	T	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT
11	[AGAT]11	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT
11	[AGAT]3 AGGT [AGAT]7	AGAT	AGAT	AGAT	AG	G	T	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT
12	[AGAT]12	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT
13	[AGAT]13	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT
14	[AGAT]14	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT	AGAT

8 alleles by length → 10 alleles by sequence

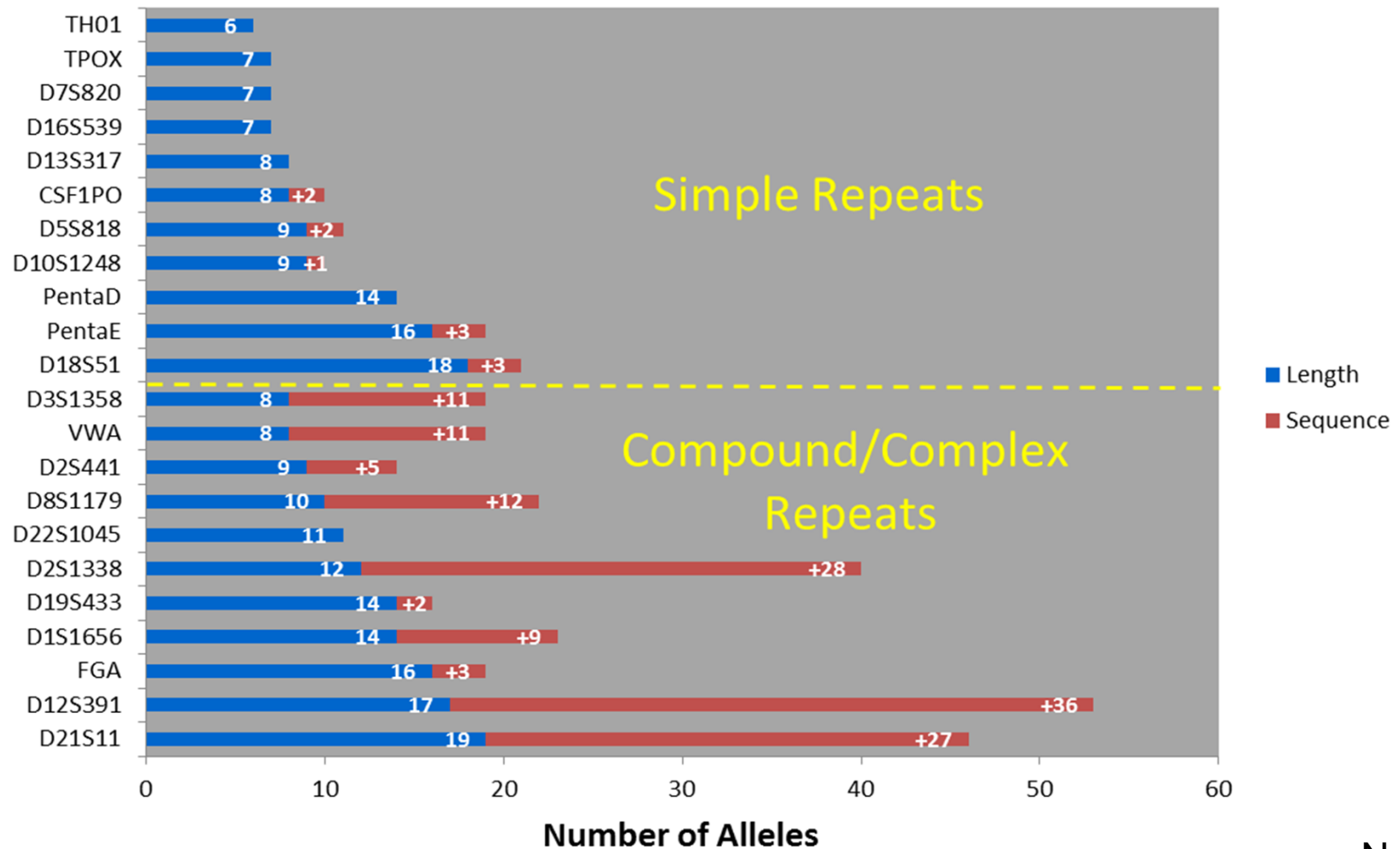
Forensic STR Sequence Diversity

Alleles Obtained by Length



Forensic STR Sequence Diversity

Alleles Obtained by Sequence

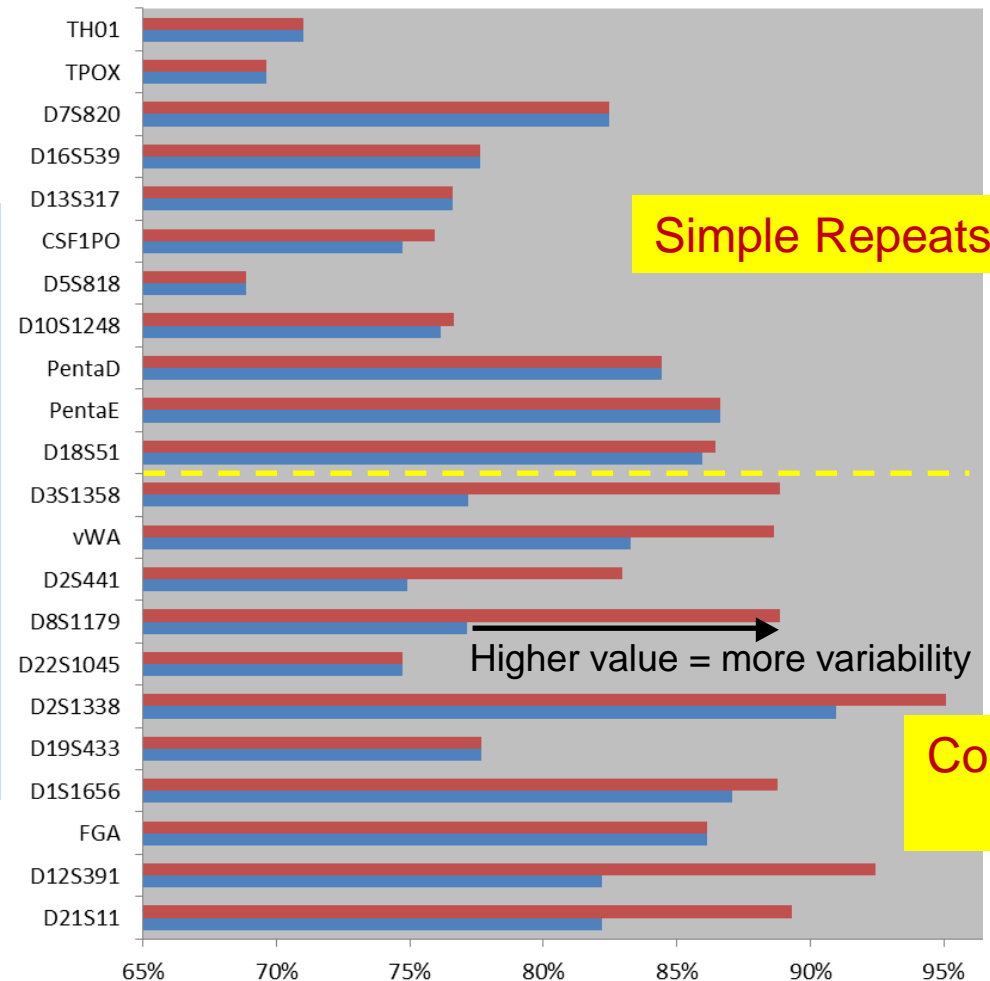


Forensic STR Sequence Diversity

Heterozygosity

$\frac{\# \text{ heterozygotes observed}}{\# \text{ of loci tested}}$

Indicates genetic
variability at a locus



Average Heterozygosity Across Populations

■ Avg Het by sequence ■ Avg Het by length

N=183

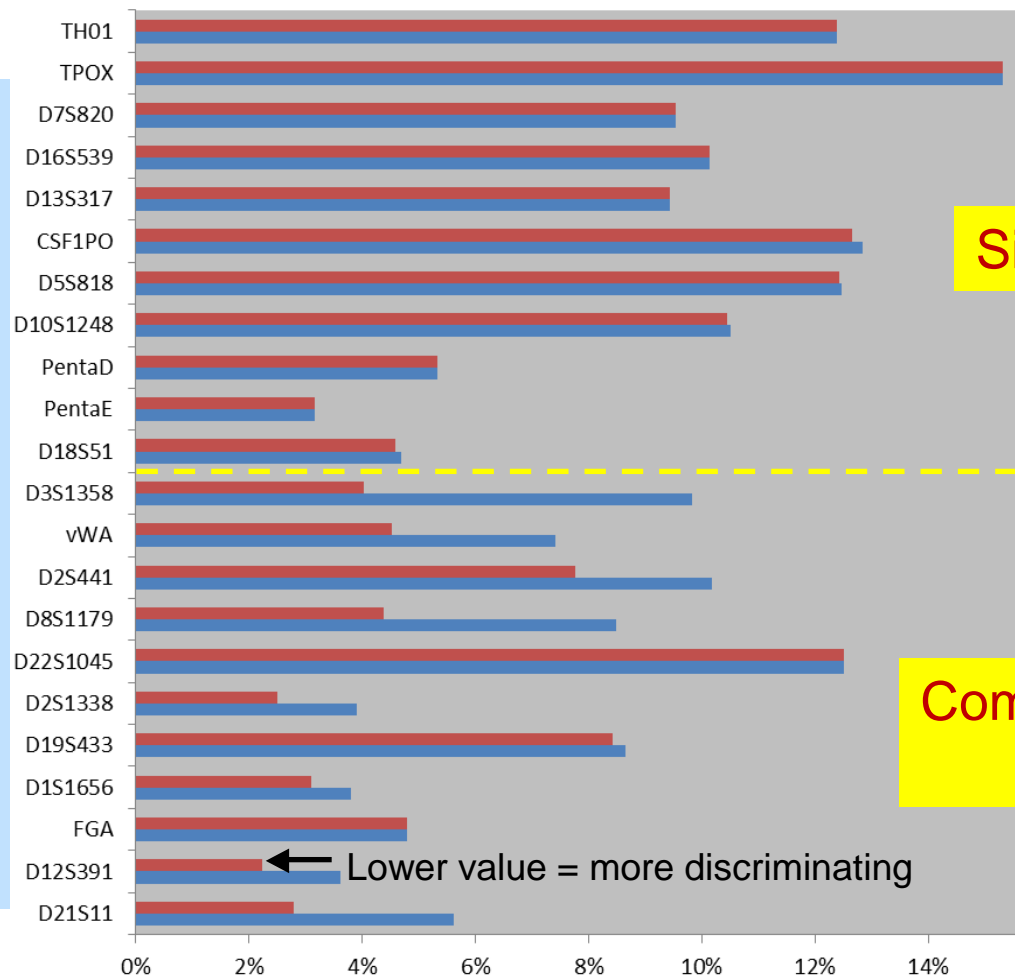
Forensic STR Sequence Diversity

Probability of Identity

Sum of each genotype frequency² at each locus

$$= \sum x_i^2$$

Probability that two unrelated individuals selected at random will have the same genotype at a locus



Simple Repeats

Compound/Complex Repeats

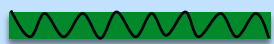
Average Probability of Identity Across Populations

■ Avg PI by sequence ■ Avg PI by length

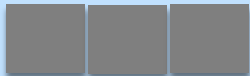
Recognition Site-Based Informatics for STRs



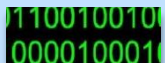
Moving recognition sites out will capture information within the flanking regions



PCR Primers



STR Repeat Region



Recognition site (~10 nt)

D7S820



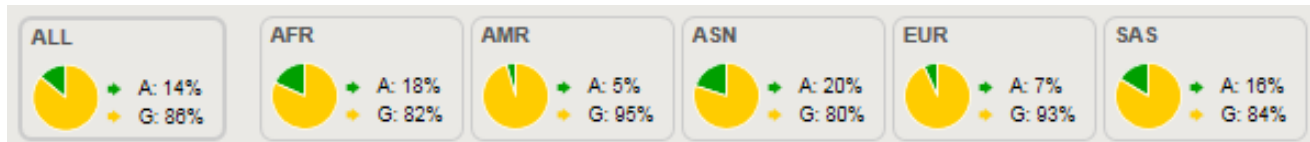
Three flanking region SNPs identified

Forensic STR Flanking Region

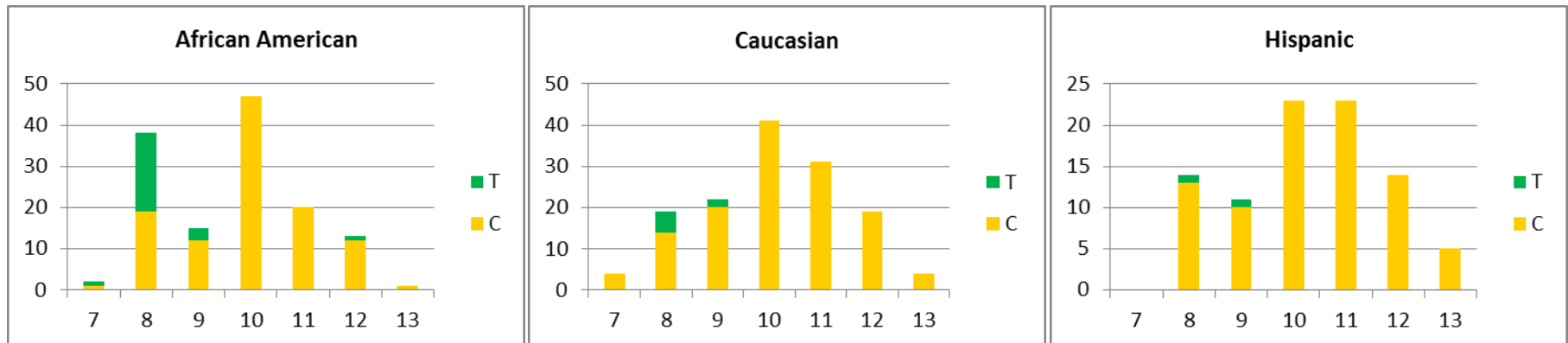
D7S820

Distribution of *rs16887642* by Population in 1000 Genomes Project

rs16887642 (note- opposite strand A/G reported):



Distribution of *rs16887642* by Population / Allele in N=183

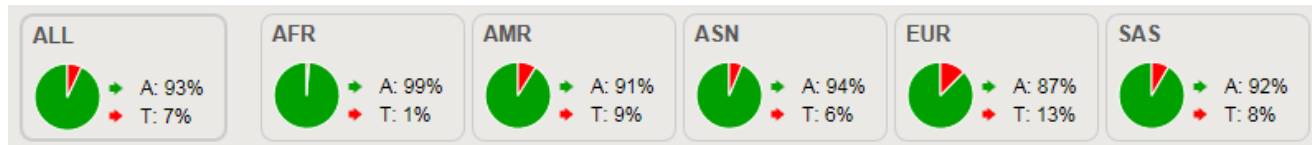


Forensic STR Flanking Region

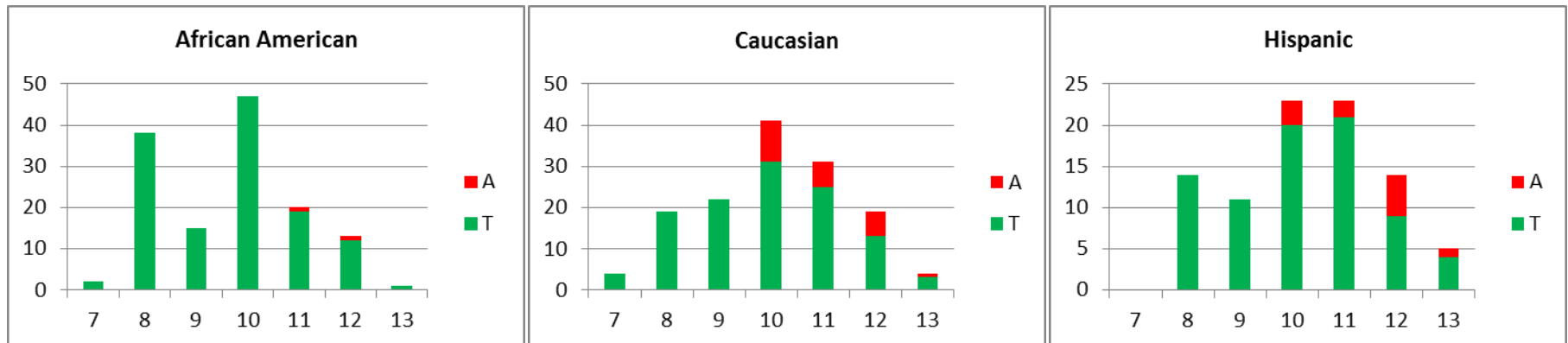
D7S820

Distribution of *rs7789995* by Population in 1000 Genomes Project

rs7789995 (note- opposite strand A/T reported):



Distribution of *rs7789995* by Population / Allele in N=183

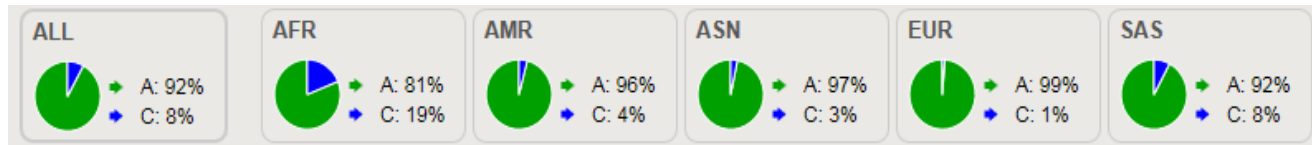


Forensic STR Flanking Region

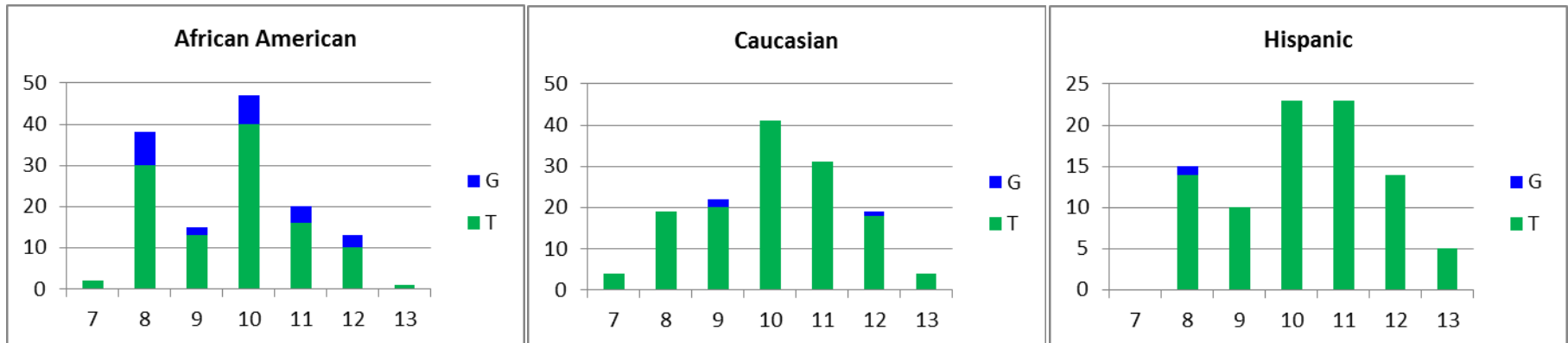
D7S820

Distribution of *rs7786079* by Population in 1000 Genomes Project

rs7786079 (note- opposite strand A/C reported):



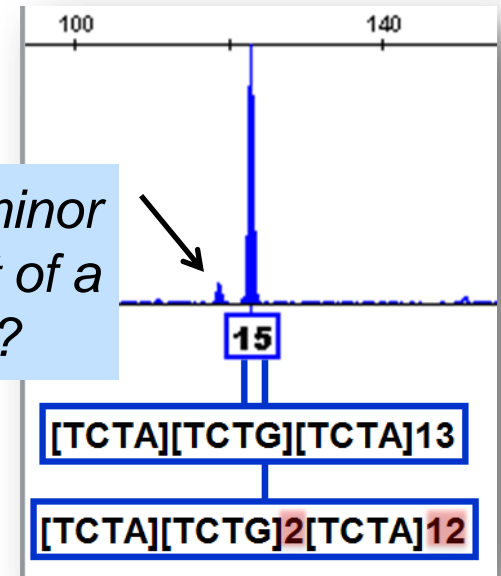
Distribution of *rs7786079* by Population / Allele in N=183



Isoalleles

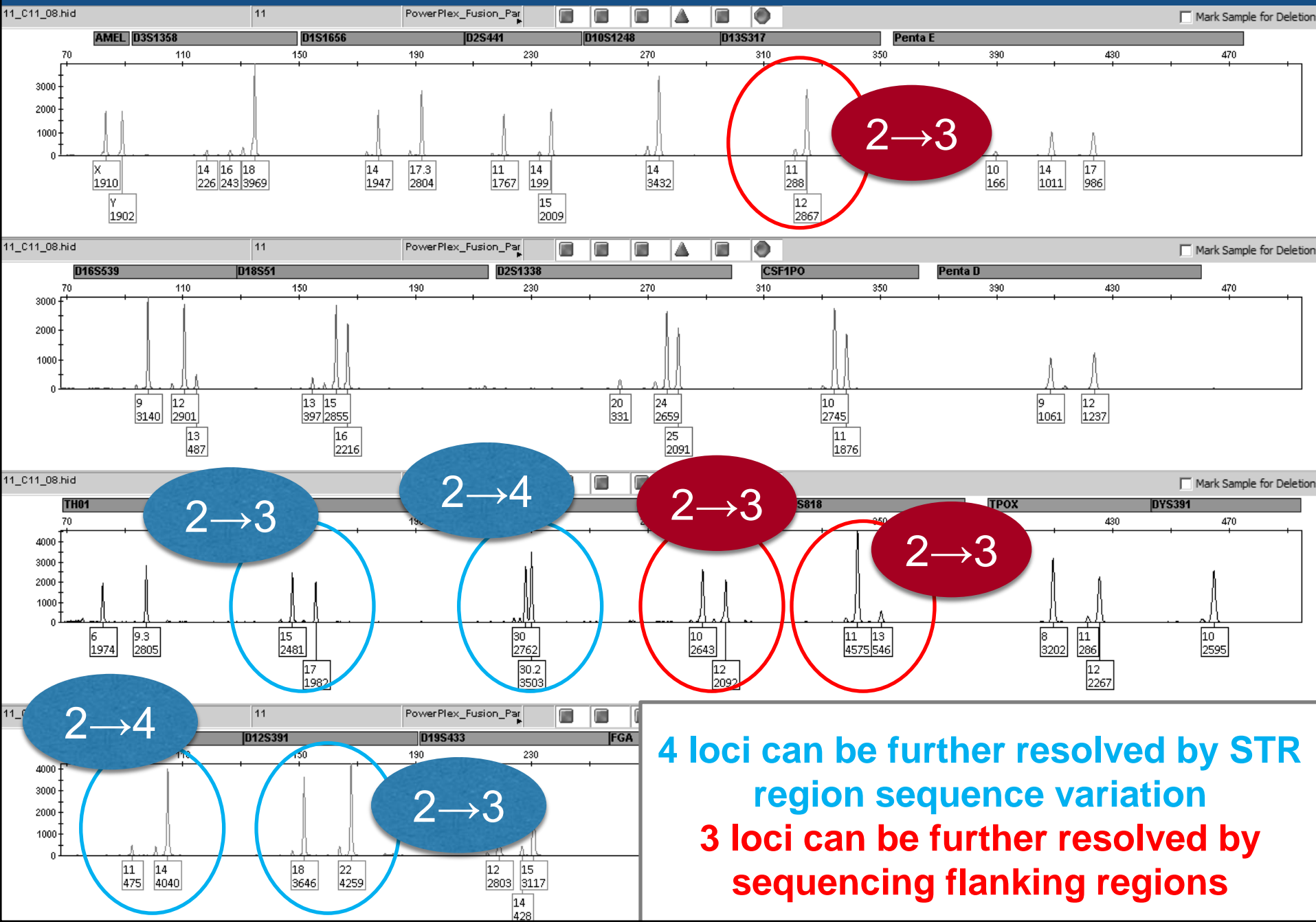
- Identical by length; unique by sequence
- Within an individual
 - Sequence based heterozygote
 - Increased heterozygosity ($p^2 \rightarrow 2pq$)
 - Marginal benefits for one-to-one matching?

Stutter or minor component of a mixture?



- Between individuals (DNA mixtures)
 - Resolve overlapping alleles 15 \rightarrow '15A' and '15B'
 - Resolve stutter from minor components

2 person, 9:1 Mixture, selected for maximal overlapping alleles



**4 loci can be further resolved by STR
region sequence variation**
**3 loci can be further resolved by
sequencing flanking regions**

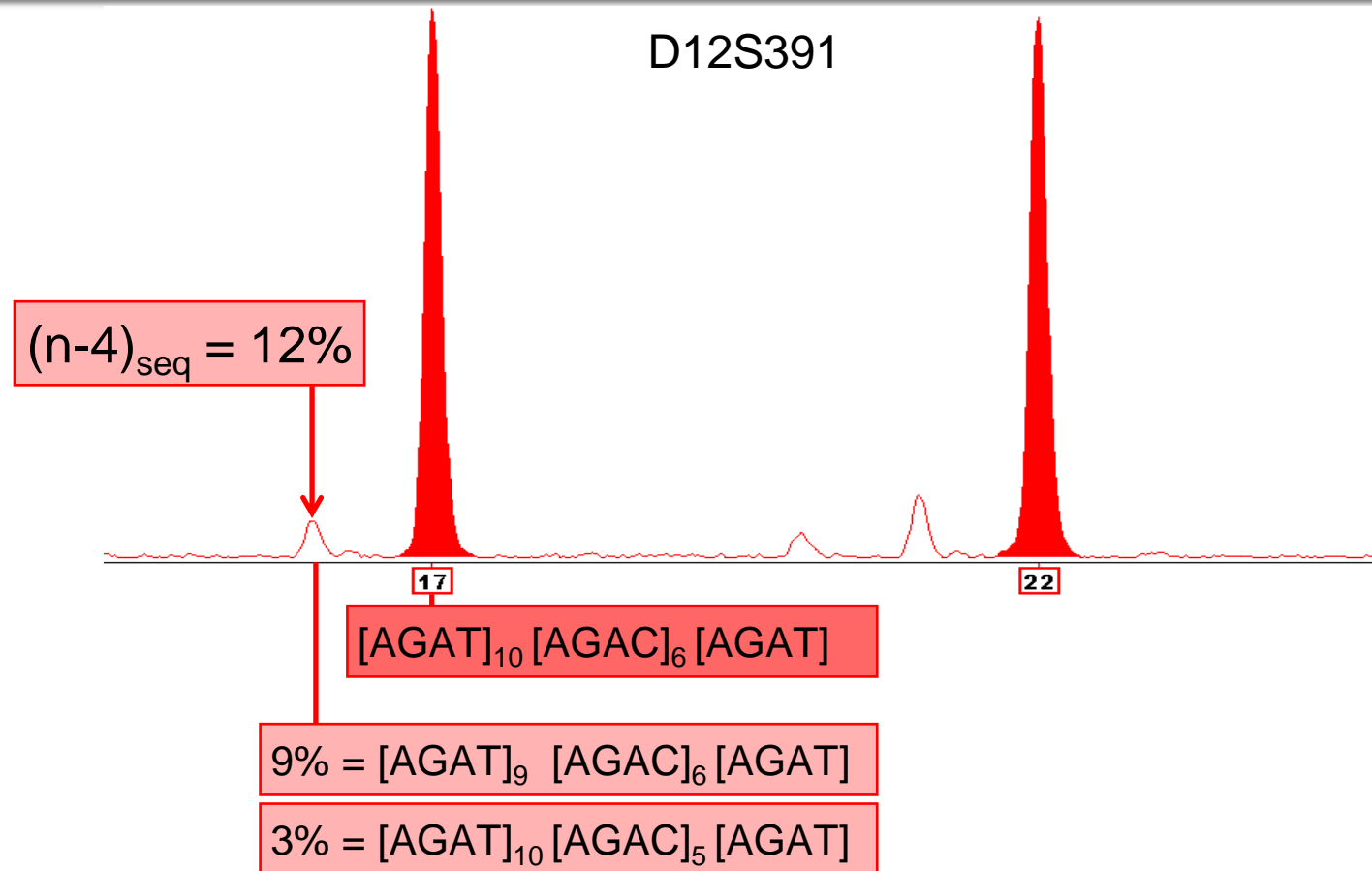
But before moving forward

In order to understand the true benefits of sequencing to mixtures

- Validate analytical and stochastic calling thresholds for NGS work
- Assess sensitivity of NGS methods to detect minor components
 - Still PCR front end – will we observe better than 10:1?
- In practice - how often will *resolvable* overlapping alleles be observed in a mixture?
- Need the allele frequencies for 'new' alleles
 - What size population databases are needed? Greater than 200?
- Incorporate NGS sequence data into probabilistic genotyping software (STRMix, True Allele, etc)
 - What are the gains in stats (Log LR)? Improved contributor ratio estimates?

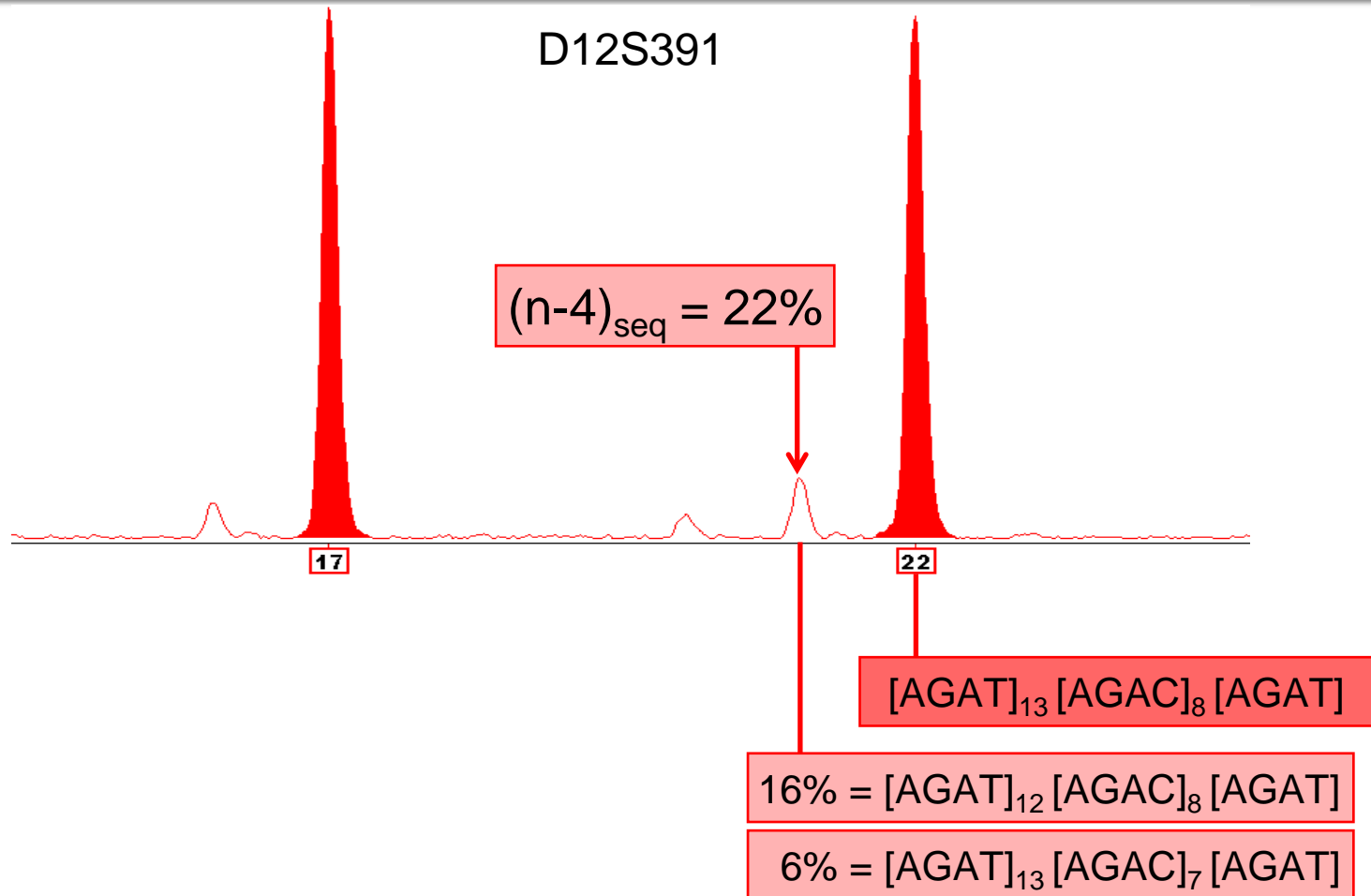
More loci – better loci?

Stutter by Sequence



Single source sample

Stutter by Sequence



Single source sample

NIST Support for NGS Research

Annotations

- GRCh38 genome with STR repeat region and flanking SNPs identified
- SNPs found in 1000 genome data set within 500 bp of STRs
- Useful for NGS primer design and bioinformatics

Excel Workbook

- Sheet for each STR locus
- Observed alleles, repeat structure, platform
- Broken out by sub-motif
- References
- Not frequency data
- This can be updated as needed

Kathe

¹National I

²G

³National I

Butler³

ryland, USA

USA

ryland, USA

Updat

ions for

NIST Support for NGS Research

SRM 2391c

SRM 2391c – Component B				
Marker	Length-based Types	Sanger Result	Repeat Structure –Allele 1	Repeat Structure –Allele 2
D1S1656	11, 14	11, 14	[TAGA] ₁₁ [TG] ₅	[TAGA] ₁₄ [TG] ₅
D2S1338	17, 17	17, 17	[TGCC] ₆ [TTCC] ₁₁	[TGCC] ₆ [TTCC] ₁₁
D2S441	10, 14	10, 14	[TCTA] ₁₀	[TCTA] ₁₁ TTTA [TCTA] ₂
D3S1358	15, 19	15, 19	TCTA [TCTG] ₃ [TCTA] ₁₁	TCTA [TCTG] ₃ [TCTA] ₁₅
D5S818	12, 13	12, 13	[AGAT] ₁₂	[AGAT] ₁₃
D6S1043	14, 19	14, 19	[AGAT] ₁₄	[AGAT] ₁₃ ACAT [AGAT] ₅
D7S820	10, 10	10, 10	[GATA] ₁₀	[GATA] ₁₀
D8S1179	10, 13	10, 13	[GATA] ₁₀	[GATA] ₁₃
D8S1115	15, 17	15, 17	[GATA] ₁₅	[GATA] ₁₇
D10S1248	13, 13	13, 13	[GATA] ₁₃	[GATA] ₁₃
D12S391	19, 24	19, 24	[GATA] ₁₉	[GATA] ₂₄
D13S317	9, 12	9, 12	[GATA] ₉	[GATA] ₁₂
D16S539	10, 13	10, 13	[GATA] ₁₀	[GATA] ₁₃
D18S51	13, 16	13, 16	[GATA] ₁₃	[GATA] ₁₆
D19S433	16, 16.2	16, 16.2	[GATA] ₁₆	[GATA] _{16.2}
D21S11	32, 32.2	32, 32.2	[GATA] ₃₂	[GATA] _{32.2}
D22S1045	15, 17	15, 17	[GATA] ₁₅	[GATA] ₁₇
CSF1PO	10, 11	10, 11	[GATA] ₁₀	[GATA] ₁₁
FGA	20, 23	20, 23	[GATA] ₂₀	[GATA] ₂₃
Penta D	8, 12	8, 12	[GATA] ₈	[GATA] ₁₂
Penta E	7, 15	7, 15	[GATA] ₇	[GATA] ₁₅
SE33	17, 18	17, 18	[GATA] ₁₇	[GATA] ₁₈
TH01	6, 9.3	6, 9.3	[GATA] ₆	[GATA] _{9.3}
TPOX	8, 11	8, 11	[GATA] ₈	[GATA] ₁₁
vWA	17, 18	17, 18	[GATA] ₁₇	[GATA] ₁₈



National Institute of Standards & Technology

Certificate of Analysis

Standard Reference Material® 2391c

PCR-Based DNA Profiling Standard

(*) Deletion of 2 bp in an uncounted repeat unit results in the 16.2 designation.

Note: Sequence information in gray indicates bases that are not counted toward the length-based genotype designation.

SRM 2391c – Component B			
Marker	Length-based Types	Sanger Result	Repeat Structure –Allele 1
DYS19	14	14	[TAGA] ₃ TAGG [TAGA] ₁₁
DYS385a	13	13	[GAAA] ₁₃
DYS385b	17	17	[GAAA] ₁₇
DYS389I	13	13	[TCTG] ₃ [TCTA] ₁₀
DYS389II	31	31	[TCTG] ₆ [TCTA] ₁₂ N ₄₈ [TCTG] ₃ [TCTA] ₁₀
DYS390	23	23	[TCTG] ₈ [TCTA] ₁₀ TCTG [TCTA] ₄
DYS391	10	10	[TCTA] ₁₀
DYS392	11	11	[TAT] ₁₁
DYS393	12	12	[AGAT] ₁₂
DYS394	13	13	[TCTA] ₈ [TCTG] ₂ [TCTA] ₄
DYS395	14	14	[TTTTC] ₁₀
DYS396	15	15	[AGAT] ₁₁
DYS397	16	16	[AGAGAT] ₁₂ N ₄₂ [AGAGAT] ₈
DYS398	17	17	[TTTC] ₁₁ N ₅₀ [TTTC] ₁₅
DYS399	18	18	[AGAT] ₁₅
DYS400	19	19	[GAAA] ₁₅ AA [GAAA] ₂
DYS401	20	20	[ATAG] ₁₀
DYS402	21	21	[CTT] ₂₅
DYS403	22	22	[G] ₃ GAAG [AAAG] ₁₄ GGAG [AAAG] ₄ N ₆
DYS404	23	23	[AAAG] ₁₅
DYS405	24	24	[ATCT] ₁₁
DYS406	25	25	[GATA] ₁₂
DYS407	26	26	[TTTC] ₁₈
DYS408	27	27	[AAAG] ₁₇
DYS409	28	28	[AGAG] ₃ [AAAG] ₁₉
DYS410	29	29	[TCTA] ₄ [TGTA] ₂ [TCTA] ₂ [TGTC] ₂ [TCTA] ₁₀
DYS411	30	30	[CTTT] ₉
DYS412	31	31	[AAAG] ₃ GTAG [GAAG] ₄ N ₂₀ [GAAG] ₉
DYS413	32	32	[AAAG] ₁₃
DYS414	33	33	[AAAG] ₃ GTAG [GAAG] ₄ N ₂₀ [GAAG] ₁₀
DYS415	34	34	[AAAG] ₁₅
Y GATA H4	11	11	[TAGA] ₁₁

Note: Sequence information in gray indicates bases that are not counted toward the length-based genotype designation.

<http://www.nist.gov/srm/>

Acknowledgements

NIST

Katherine Gettings

Nate Olson

Jo Lynne Harenza

Mike Coble

Becky Steffen

Margaret Kline

Student Interns

Rachel Aponte (GWU)

Harish Swaminathan

(Rutgers)

Anna Blendermann (MC)

Funding

**FBI – DNA as a
Biometric**

Battelle

Seth Faith (now @NCSU)

Rich Guerrieri

Brian Young

Liz Montano

Esley Heizer

Angela Minard-Smith

Christine Baker

Promega

Doug Storts

Jay Patel

Contact Information

peter.vallone@nist.gov

